# Subspace Models
# for Acoustic Unit Discovery

Lucas Ondel, Bolaji Yusuf

Brno University of Technology, Faculty of Information Technology
Božetěchova 1/2. 612 66 Brno - Královo Pole

{iondel,xyusuf00}@fit.vutbr.cz

BRNO FACULTY
UNIVERSITY OF INFORMATION
OF TECHNOLOGY TECHNOLOGY

December 10, 2020

# Table of Contents

- Audio recordings without labels
- Inventory of phone-like units (we call them "acoustic units")
- Segmentation and labelling

# Documenting endangered languages

- Language diversity is diminishing world wide
  - Speech technologies are only available for a few languages
  - Majority of languages are not written: main-stream ASR is not applicable
- an accurate AUD system could:
  - help linguists to document languages
  - serve as a front-end of speech-technologies for non-written languages
- 2022-2032: The decade of Indigenous Languages

  ▸ UNESCO website

# Computational model human learning

- The learning cognitive process of humans remains largely unknown:
    - the brain is very complex
    - "sensory learning" happens at a very early age, when the infants cannot communicate verbally
- Reverse engineering approach ("E. Dupoux. 2018") ▸ Link :
    - let's build model to learn speech in an unsupervised fashion
    - analyse it and derive the learning principle
    - pave the way to a more realistic artificial intelligence

- The "always more data" approach raises concerns:
  - social/ethical problems: monopoly of ML technologies by LARGE data owner
  - ecological issues: more data implies more energy consumption
- AUD research implies data efficient models



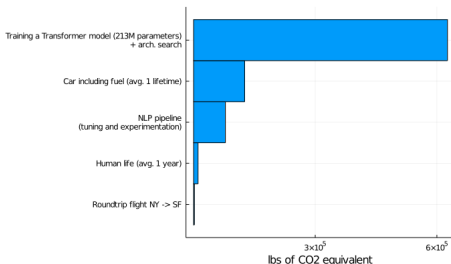Figure: Energy consumption of training deep learning models. Source: Strubell *et al*, 2019 ▸ Link

# Outline

- heuristic-based model: $\sim 1990 - 2005$
- non-parametric Bayesian-based models: $\sim 2005 - 2020$
- neural network-based models: $\sim 2015 - 2020$

- VAE-HMM: Variational AutoEncoder with HMM/GMM prior over the latent space. ▸ Link
- Visually guided neural network: replace textual transcription by images ▸ Link
- VQ-VAE: AutoEncoder with discretized bottleneck ▸ Link

# Non-parametric Bayesian-based models

- Non-parametric model for word segmentation  ▸ Link 1  ▸ Link 2
- Dirichlet-Process HMM (2012)  ▸ Link
- Dirichlet-Process HMM (Variational Bayes inference) (2016)
  ▸ Link

| p1 | p2 | p3 |
|---|---|---|
| au1 | | au2 |

- Clustering:
  - Normalized Mutual Information: $200 \frac{MI(X,Y)}{H(X)+H(Y)} \%$
  - Measures the statistical relation between the discovered units and the ground truth transcription
  - 100 % → one-to-one mapping between AU and phones
  - 0 % → AUs are completely uninformative
- Segmentation:
  - F-score: harmonic mean of segmentation precision and recall
  - ± 20 ms tolerance allowed

- Data:
  - Mboshi (Congo Brazzaville) - 3-4 hours
  - Yoruba (West Africa - Nigeria) - 10 hours
  - English (TIMIT) - 4 hours
- Features:
  - traditional features: MFCCs + Δ + ΔΔ

- Adapted version of "Chorowski et al., 2019" ▸ Link
- Dirichlet Process HMM (VB inference) ▸ Link

| Corpus | System | NMI | F-Score | # units |
|--------|--------|-----|---------|---------|
| English | VQ-VAE | 29.73 | 38.59 | 100 |
| English | HMM | **35.47** | **63.03** | 95 |
| Mboshi | VQ-VAE 64 | 26.85 | 20.22 | 100 |
| Mboshi | HMM | **36.47** | **47.93** | 94 |
| Yoruba | VQ-VAE 64 | 29.36 | 7.74 | 100 |
| Yoruba | HMM | **36.71** | **28.47** | 95 |

Table: Comparison of the HMM vs the VQ-VAE baseline

# Example

HMM



VQ-VAE

# Results

- VQ-WAV2VEC "A. Baevski et al., 2020" trained on 960h of LibriSpeech (unsupervised) ▸ Link
- Dirichlet Process HMM (VB inference) ▸ Link

| Corpus | System | NMI | F-Score | # units |
|--------|--------|-----|---------|---------|
| English | VQ-WAV2VEC (Gumbel) | 35.20 | 26.84 | 12008 |
| English | VQ-WAV2VEC (K-mean) | 34.06 | 25.64 | 20057 |
| English | HMM | **35.47** | **63.03** | 95 |

Table: Comparison of the HMM vs the VQ-VAE baseline

# Example



HMM

VQ-WAV2VEC

- Infants do not learn from scratch (Kuhl *et al*, 1992 ▶ Link ):
  - They have some innate sensitivity to human languages
  - With time, they become specialized to their native language
- Hypothesis/Design choice:
  - This innate sensitivity guide infants to learn the structure of speech
  - The AUD system should adapt and become language specific
- Proposal: we will use Bayesian Subspace Model techniques to implement these properties:
  - Subspace Hidden Markov Model (SHMM)
  - Hierarchical Subspace Hidden Markov Model (HSHMM)

$$p(\boldsymbol{\eta}|\mathbf{X}) = \frac{p(\mathbf{X}|\boldsymbol{\eta})p(\boldsymbol{\eta})}{p(\mathbf{X})} \tag{1}$$

- We want to design an educated prior over the AU's parameters : $p(\boldsymbol{\eta})$

- 

$$\mathbf{h} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \tag{2}$$

$$\boldsymbol{\eta} = f(\mathbf{W}\mathbf{h} + \mathbf{b}) \tag{3}$$

- We estimate the subspace parameters $\mathbf{W}, \mathbf{b}$ on annotated corpora
- Model is trained by optimizing the Evidence Lower-BOund (ELBO):

$$
\begin{aligned}
\mathcal{L} = &\langle \ln p(\mathbf{X}, |\mathbf{z}, \mathbf{W}, \mathbf{b}, \mathbf{h}_{1:T}) \rangle_q \\
&- D_{KL}(q(\mathbf{z}) \| p(\mathbf{z})) \\
&- D_{KL}(q(\mathbf{W})q(\mathbf{b})) \| p(\mathbf{W})p(\mathbf{b})) \\
&- D_{KL}(q(\mathbf{h}_{1:T}) \| p(\mathbf{h}_{1:T}))
\end{aligned}
$$

- The training follows an Expectation-Maximization-like training:
  - E-step: Baum-Welch algorithm to estimate states' occupancy
  - M-step: No closed form solution, using re-parameterization trick.

- The subspace parameters $\mathbf{W}, \mathbf{b}$ are fixed, we just learn the embeddings $\mathbf{h}$ on the target language
- Model is trained by optimizing the Evidence Lower-BOund (ELBO):

$$
\begin{aligned}
\mathcal{L} = &\langle \ln p(\mathbf{X}, |\mathbf{z}, \mathbf{W}, \mathbf{b}, \mathbf{h}_{1:T}) \rangle_q \\
&- D_{KL}(q(\mathbf{z})||p(\mathbf{z})) \\
&- D_{KL}(q(\mathbf{h}_{1:T})||p(\mathbf{h}_{1:T}))
\end{aligned}
$$

- The training follows an Expectation-Maximization-like training:
  - E-step: Baum-Welch algorithm to estimate states' occupancy
  - M-step: No closed form solution, using re-parameterization trick.

▶ Link

- Data:
    - Source languages (transcribed)
        - French, German, Polish, Spanish from Globalphone
        - 3-4 hours subsets of each language's training data
    - Target languages (untranscribed)
        - Mboshi (Congo Brazzaville) - 3-4 hours
        - Yoruba (West Africa - Nigeria) - 10 hours
        - English (TIMIT) - 4 hours
- Features:
    - traditional features: MFCCs + $\Delta$ + $\Delta\Delta$

| Corpus | System | Training | NMI | F-Score |
|--------|--------|----------|-------|---------|
| English | HMM | no | 1.74 | 0.20 |
| English | SHMM | no | 20.83 | 58.94 |
| Mboshi | HMM | no | 1.65 | 0.02 |
| Mboshi | SHMM | no | 21.0 | 39.28 |
| Yoruba | HMM | no | 4.43 | 1.26 |
| Yoruba | SHMM | no | 26.1 | 27.45 |

Table: Comparison of the HMM vs the SHMM before training

| Corpus | System | Training | NMI | F-Score |
|--------|--------|----------|--------|---------|
| English | HMM | no | 1.74 | 0.20 |
| English | HMM | yes | 35.47 | 63.03 |
| English | SHMM | no | 20.83 | 58.94 |
| English | SHMM | yes | **39.66** | **75.92** |
| Mboshi | HMM | no | 1.65 | 0.02 |
| Mboshi | HMM | yes | 36.47 | 47.93 |
| Mboshi | SHMM | no | 21.0 | 39.28 |
| Mboshi | SHMM | yes | **38.42** | **57.26** |
| Yoruba | HMM | no | 4.43 | 1.26 |
| Yoruba | HMM | yes | 35.27 | 28.83 |
| Yoruba | SHMM | no | 26.1 | 27.45 |
| Yoruba | SHMM | yes | **37.56** | **36.64** |

Table: Comparison of the HMM vs the SHMM before training

- We assume the subspace is known and fixed during AUD
- Subspace is the same for all the target languages

- We would like to adapt the subspace to make it language specific

- We would like to adapt the subspace to make it language specific

- We would like to adapt the subspace to make it language specific

- We would like to adapt the subspace to make it language specific

- We want to design an educated prior over all possible subspace: $p(\mathbf{W}, \mathbf{b})$

-

$$\boldsymbol{\alpha} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \tag{4}$$

$$\mathbf{W} = \sum_{i=1}^{Q} \alpha_i \mathbf{M}_i \tag{5}$$

$$\mathbf{b} = \sum_{i=1}^{Q} \alpha_i \mathbf{m}_i \tag{6}$$

$$\mathbf{h} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \tag{7}$$

$$\boldsymbol{\eta} = f(\mathbf{W}\mathbf{h} + \mathbf{b}) \tag{8}$$

$$\boldsymbol{\alpha} = \left[\alpha_1, \alpha_2\right]^\top \quad (9)$$
$$\mathbf{W} = \alpha_1\mathbf{M}_1 + \alpha_2\mathbf{M}_2 \quad (10)$$

- We estimate the "hyper-subspace" parameters $\mathbf{M}, \mathbf{m}, \boldsymbol{\alpha}$ on annotated corpora

- Model is trained by optimizing the Evidence Lower-BOund (ELBO):

$$\mathcal{L} = \langle \ln p(\mathbf{X}, |\mathbf{z}, \mathbf{M}_{1:Q}, \mathbf{m}_{1:Q}, \mathbf{h}_{1:T}, \boldsymbol{\alpha}) \rangle_q$$
$$- D_{\mathsf{KL}}(q(\mathbf{z})||p(\mathbf{z}))$$
$$- D_{\mathsf{KL}}(q(\mathbf{M}_{1:Q})q(\mathbf{m}_{1:Q})||p(\mathbf{M}_{1:Q})p(\mathbf{m}_{1:Q}))$$
$$- D_{\mathsf{KL}}(q(\mathbf{h}_{1:T})||p(\mathbf{h}_{1:T}))$$
$$- D_{\mathsf{KL}}(q(\boldsymbol{\alpha})||p(\boldsymbol{\alpha}))$$

- The training follows an Expectation-Maximization-like training:
  - E-step: Baum-Welch algorithm to estimate states' occupancy
  - M-step: No closed form solution, using re-parameterization trick.

- The "hyper-subspace" parameters $\mathbf{M}, \mathbf{m}$ are fixed, we just learn the embeddings $\mathbf{h}$ and $\alpha$ on the target language
- Model is trained by optimizing the Evidence Lower-BOund (ELBO):

$$
\begin{aligned}
\mathcal{L} = & \langle \ln p(\mathbf{X}, |\mathbf{z}, \mathbf{M}_{1:Q}, \mathbf{m}_{1:Q}, \mathbf{h}_{1:T}, \alpha) \rangle_q \\
& - D_{KL}(q(\mathbf{z}) || p(\mathbf{z})) \\
& - D_{KL}(q(\mathbf{h}_{1:T}) || p(\mathbf{h}_{1:T})) \\
& - D_{KL}(q(\alpha) || p(\alpha))
\end{aligned}
$$

- The training follows an Expectation-Maximization-like training:
  - E-step: Baum-Welch algorithm to estimate states' occupancy
  - M-step: No closed form solution, using re-parameterization trick.

▶ Link

- Data:
  - Source languages (transcribed)
    - French, German, Polish, Spanish from Globalphone
    - 3-4 hours subsets of each language's training data
  - Target languages (untranscribed)
    - Mboshi (Congo Brazzaville) - 3-4 hours
    - Yoruba (West Africa - Nigeria) - 10 hours
    - English (TIMIT) - 4 hours
- Features:
  - traditional features: MFCCs + $\Delta$ + $\Delta\Delta$

| Corpus | System | NMI | F-Score |
|--------|--------|-------|---------|
| English | HMM | 35.47 | 63.03 |
| English | SHMM | 39.66 | 75.92 |
| Mboshi | HMM | 36.47 | 47.93 |
| Mboshi | SHMM | 38.42 | 57.26 |
| Yoruba | HMM | 35.27 | 28.83 |
| Yoruba | SHMM | 37.56 | 36.64 |

Table: Comparison of the HMM, SHMM and HSHMM

| Corpus | System | NMI | F-Score |
|--------|--------|-------|---------|
| English | HMM | 35.47 | 63.03 |
| English | SHMM | 39.66 | 75.92 |
| English | SHMM (2) | 37.46 | 72.19 |
| | | | |
| Mboshi | HMM | 36.47 | 47.93 |
| Mboshi | SHMM | 38.42 | 57.26 |
| Mboshi | SHMM (2) | 35.50 | 51.28 |
| | | | |
| Yoruba | HMM | 35.27 | 28.83 |
| Yoruba | SHMM | 37.56 | 36.64 |
| Yoruba | SHMM (2) | 35.72 | 31.34 |

Table: Comparison of the HMM, SHMM and HSHMM

SHMM (2): the subspace is retrained on the target language.

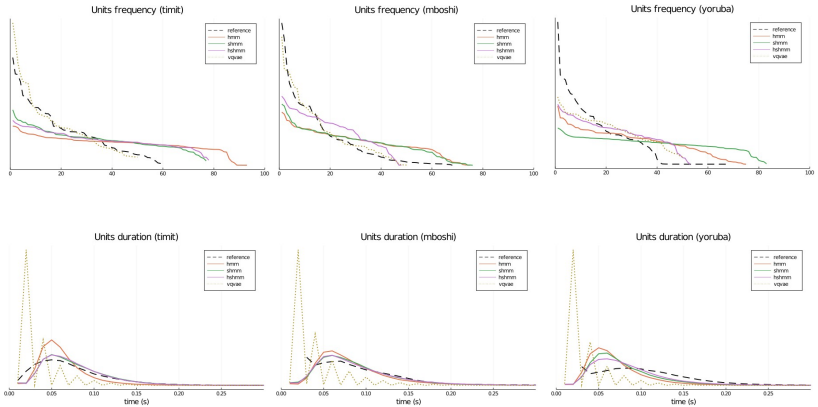| Corpus | System | NMI | F-Score |
|---|---|---|---|
| English | HMM | 35.47 | 63.03 |
| English | SHMM | 39.66 | 75.92 |
| English | SHMM (2) | 37.46 | 72.19 |
| English | HSHMM | **40.56** | **78.32** |
| Mboshi | HMM | 36.47 | 47.93 |
| Mboshi | SHMM | 38.42 | 57.26 |
| Mboshi | SHMM (2) | 35.50 | 51.28 |
| Mboshi | HSHMM | **41.17** | **60.82** |
| Yoruba | HMM | 35.27 | 28.83 |
| Yoruba | SHMM | 37.56 | 36.64 |
| Yoruba | SHMM (2) | 35.72 | 31.34 |
| Yoruba | HSHMM | **37.88** | **38.44** |

Table: Comparison of the HMM, SHMM and HSHMM

SHMM (2): the subspace is retrained on the target language.

# Some statistics

HMM

SHMM

HSHMM

- We have proposed two new models for the task of Acoustic Unit Discovery:
  - Subspace Hidden Markov Model
  - Hierarchical Subspace Hidden Markov Model
- These models are inspired by how infants learn to speak
- They show strong improvement in terms of clustering and segmentation quality
- The concept of (hierarchical) subspace and can be extended to a large class of models
- To reproduce our experiments:
  ```
  https://github.com/beer-asr
  ```

- AUD is not a solved problem!
- Model suffers from high variability of speech
- Two major problems:
  - Acoustic modeling: going beyond HMM
  - Language modeling: discovery words
- Towards the first system to learn speech as humans...

Lucas Ondel, Hari Krishna Vydana, Lukáš Burget, and Jan Černocký (2019). "Bayesian Subspace Hidden Markov Model for Acoustic Unit Discovery". In: Proc. Interspeech 2019, pp. 261–265. URL: http://dx.doi.org/10.21437/Interspeech.2019-2224.

Lucas Ondel, Pierre Godard, Laurent Besacier, Elin Larsen, Mark Hasegawa-Johnson, Odette Scharenborg, Emmanuel Dupoux, Lukáš Burget, Francois Yvon, and Sanjeev Khudanpur (2018). "Bayesian Models for Unit Discovery on a Very Low Resource Language". In: 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, pp. 5939–5943. URL: https://ieeexplore.ieee.org/document/8461545.

Lucas Ondel, Lukaš Burget, Santosh Kesiraju, and Jan Černocký (2017). "Bayesian phonotactic language model for acoustic unit discovery". In: 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, pp. 5750–5754. URL: `https://www.fit.vut.cz/research/publication/11472/.en`.

Lucas Ondel, Lukáš Burget and Jan Černocký (2016). "Variational inference for acoustic unit discovery". In: Procedia Computer Science 81, pp. 80–86. URL: `https://www.sciencedirect.com/science/article/pii/S1877050916300473`.

Thank you for your attention.